A Future for Learning Semantic Models of Man-Made Environments

Wolfgang Förstner

University of Bonn, Department of Geodesy and Geoinformation D-53121 Bonn, Nussallee 15, Email: wf@ipb.uni-bonn.de

Abstract-Deriving semantic 3D models of man-made environments hitherto has not reached the desired maturity which makes human interaction obsolete. Man-made environments play a central role in navigation, city planning, building management systems, disaster management or augmented reality. They are characterised by rich geometric and semantic structures. These cause conceptual problems when learning generic models or when developing automatic acquisition systems. The problems appear to be caused by (1) the incoherence of the models for signal analysis, (2) the type of interplay between discrete and continuous geometric representations, (3) the inefficiency of the interaction between crisp models, such as partonomies and taxonomies, and soft models, mostly having a probabilistic nature, and (4) the vagueness of the used notions in the envisaged application domains. The paper wants to encourage the development and learning of generative models, specifically for man-made objects, to be able to understand, reason about, and explain interpretations.

I. INTRODUCTION

Deriving semantic 3D models of man-made environments has gained interest since the beginning of image analysis, see (Brooks, 1983; Herman and Kanade, 1986) and the surveys for outdoor and indoor environments in (Musialski et al., 2013; Fidler and Urtasun, 2015). Man-made environments play a central role in navigation, city planning, building management systems, disaster management or augmented reality.

Automatic methods for semantic building reconstruction hitherto have not reached the desired maturity which makes human interaction obsolete. In spite of great success in automatically reconstructing the geometry of buildings it appears that the rich geometric and semantic structures, which characterize man-made objects, slows down progress. The paper identifies successes and difficulties in using explicit models for supporting the geometric and semantic reconstruction of buildings. We want to encourage the development and learning of generative models, specifically for man-made objects, be able to understand, reason about, and explain interpretations of man-made scenes, quite in the spirit of (Lake et al., 2016).

Based on experiences in our research group, we will discuss typical tasks which we solved using structural descriptions (image orientation, building reconstruction, and façade interpretation), and embed the used methods in the stream of concurrent solutions. We try to identify the attempts to learn the underlying models and the achievements in object recognition which on one had promise to support future methods for interpreting images of man-made scenes. However, these models – in our view – still contain conceptual problems when learning generative models or when developing automatic acquisition systems. The problems appear to be caused by (1) the incoherence of the models for signal analysis, (2) the type of interplay between discrete and continuous geometric representations, (3) the inefficiency of the interaction between crisp models, such as partonomies and taxonomies, and soft models, mostly having a probabilistic nature, and (4) the vagueness of the used notions in the envisaged application domains. A goal for future research should be to learn building models, i.e., to learn geometric, structural and semantic models which help understanding images of man-made scenes, and to further develop methods to learn these highly structured models.

We start with experiences with structural descriptions for solving tasks related to man-made objects.

II. USING STRUCTURAL DESCRIPTIONS

In the following we discuss three basic problems, pose determination, building reconstruction, and image interpretation in the context of man-made scenes. Pose determination is representative for the large class of parameter estimation problems based on correspondences, where - depending on the number of available image features - structural descriptions may be of advantage. Building reconstruction is a representative for the large class of problems where, besides a large number of parameters, also the structure of the solution, especially the number of parameters and possibly the constraints between the parameters, is not known from the beginning. Finally, image interpretation aims at a semantic description, thus above parameters and structure also aims at finding the class memberships of the objects and possibly the semantic relations between the objects shown in the images. The discussion of these tasks is triggered by own research and the solutions known before and achieved later and gives insight into the development during the last three decades w.r.t. the used representations and reasoning methods.

A. Relational Matching using Edges for Pose Determination

Pose estimation requires correspondences between images and a 3D model, which, when performed automatically, requires adequate matching techniques. Matching a given model with an image is based on a common representation. Thirty years ago, due to limited computer power, representations based on point or line type features dominated. Keypoints were mainly used for image-to-image matching whereas model-toimage matching mainly use image edges, even only straight edge segments, see e.g., (Brooks et al., 1979; Lowe, 1987) and Fig. 1. The search for correspondences was done incre-



Fig. 1. Pose estimation based on straight line segments. Left Image edges. Mid: One of the models given as 3D line segments. Right: Image with projected model; from (Lowe, 1987). The matching is based on triplets of corresponding lines, which allow to directly derive the pose parameters, which then are checked for consistency with the other edges

mentally, formalized as interpretation tree by (Grimson and Lozano-Perez, 1987).

Based on work on the consistent labelling problem (Haralick and Shapiro, 1979, 1980) and relation matching (Shapiro and Haralick, 1987), we in the late nineteen eighties explored model-to-image matching for finding buildings (roofs) in aerial images (Förstner, 1988; Schickler, 1992), and more general relational descriptions for pose estimation (Vosselman, 1992) or map-to-image matching based on road networks (Haala and Vosselman, 1992; Vosselman and Haala, 1992), see Fig. 2.

Given a wire frame model of the object and line segments together with their mutual relations, such as connectivity or parallelity, the task was to derive the six parameters of the pose. Matching of the two relational descriptions using heuristic search (A*) was based on a probabilistic model of the projection. The matching costs were based on the mutual self-information $I(x; y) = -\log(P(x)/P(x|y))$.¹ The goal of the search was to maximize the sum of the mutual self-information of all matches and relations. The probabilities were learned from training data. This simplifies the evaluation of missing correspondences – often called *wild cards* in matching – by setting I(x; y) = 0, since the missing match has no influence. An example for detecting a road junction in an aerial image is given in Fig. 2



Fig. 2. Image-to-model matching. Left: Road map as planar graph. Mid: Search tree for image orientation. **Right:** Match with aerial image; from (Vosselman and Haala, 1992). The model has 25 units (junctions, edges) (only the region around the road intersection), the image has 21 units, the search tree has 52 nodes, determining the orientation was tried six times, the software was written in POP-11, the computing time was 227 seconds on a VAX 3200

¹The mutual self-information $I(x; y) \in (-\infty, \infty)$ depends on the probabilities. The mutual entropy $H(x; y) = \mathbb{E}_{p(x,y)}(I(x; y)) \geq 0$ is its expectation and often called mutual information.

Progress in pose estimation is based on more informative features (Brachmann et al., 2016) or first estimating viewpoints using a regression convolutional neural network and then using key points for fine matching (Tulsiani and Malik, 2015). While both directions do not use an explicit model of the scene, exploiting a hierarchical object model for efficient detection (Mottaghi et al., 2015), see Fig. 3 and generalizations to articulated objects are indispensable for locating persons in general pose, see e.g., (Kar et al., 2014). Progress is triggered by a 3D recognition challenge (Xiang et al., 2014).



Fig. 3. Hierarchical model for object detection, including a step for determining the orientation of the object. **Left:** Hierarchical model with three layers. **Mid:** Given image. **Right:** Bounding box, class, and projected coarse model; from (Mottaghi et al., 2015)

B. Generic Building Models from Multiple Images using Constraint Programming

Reconstructing generic building models from images requires an adequate representation of the structure. Structure refers to the number of building parts, their relations w.r.t. neighbourhood and geometry, and to constraints between the parts, especially among parameters of the individual parts. In a first step the reconstruction only aims at a rich geometric description, and does not include an interpretation. This may be fruitful in a later step, see (Malik et al., 2015).

Early work (Brooks et al., 1979) fixed the structure and only allowed variations for parameters for parameter. The first work assuming buildings to be represented as polyhedra is (Herman and Kanade, 1986; Huertas and Nevatia, 1988). Explicitly deriving neighbourhood relations between building parts was addressed by (Fua and Hanson, 1987), see Fig. 4.



Fig. 4. Deriving the topology of a complex building. Left: Aerial image. Mid: Result of data driven segmentation. **Right:** Automatically derived symbolic image description; from (Fua and Hanson, 1987)

Based on these stimulating results and motivated by the Avenches building extraction benchmark (Mason et al., 1994) we addressed the reconstruction of complex buildings from multiple images, see (Braun et al., 1995; Fischer et al., 1998). Buildings are assumed to be hierarchically decomposed, to consist of building parts and its projection into the image yield corners, each being an aggregates of a point and its neighbouring edges and faces. The building parts are parametrized wire



Fig. 5. 3D reconstruction of complex buildings from multiple images. Left: Four image sections. Mid left: Reconstructed 3D corners. Mid right: Triggered building parts: five terminals, one connector with three faces (belonging to the blue junction). Right: Reconstructed building (roof) fulfilling topological and geometrical constraints; from (Fischer et al., 1998)

frames. The reconstruction method employees the integration of a data-driven trigger phase and a model-driven verification phase. In a first step, mutually oriented images 3D vertices were reconstructed (see Fig. 5), based on keypoints and neighbouring 3D edges in an prespecified area of interest, see (Fuchs and Förstner, 1995; Lang and Förstner, 1996). Based on the 3D vertices, building parts were hypothesized and mutual topological and geometrical constraints were exploited to reconstruct the complete building. The method was implemented as constraint satisfaction program (in constraint logic programming, see (Kolbe et al., 2000; Fischer, 2000)), and allowed for occlusions (wild cards, see above) and incrementally for the prediction of new corners. The verification step included a prior on the different buildings and viewing directions, which exploited the shortest coding of the expected feature adjacency graph, see (Heuel and Kolbe, 2001; Kolbe, 2000).

Generating highly structured city models requires a quite generic building model, with a variable number of parts. The models we used are limited. They use restricted prespecified parametrized building parts, and thus cannot be used for larger areas. Though constraint logic programming appeared to be useful, the statistical knowledge only influenced the heuristics of the search and in a prespecified manner was used in the final evaluation. The parts need to be learned together with their relations and the reconstruction should exploit the learned statistics: Neither was the likelihood of the extracted features exploited as e.g., in (Arbelaez et al., 2011), nor was any knowledge about illumination (sun angle, albedo) used. This would allow an integration of forward and backward modelling using computer graphics, see the example in Fig. 3.

There is a dichotomy: whether it is more favourable to aim at less simple parts with complex relations or to try to find more expressive complex parts with more restricted relations, with the inherent question how to deal with curved surfaces then. The extreme, representing the surface as mesh up to now appears to be the most flexible and successful approach. This circumvents the problem of structuring, which then needs to be addressed in a second step, see the next section.

Later work on building reconstruction exploited regularities of roof tops based on the straight skeleton (Brenner, 2000) or aimed at watertight reconstruction for outdoor (Zhou and Neumann, 2010) or indoor scenes from LiDAR data (Oesau et al., 2013). We can observe intensive research in reconstructing large city areas based on terrestrial and aerial images using classical pipelines, which are developed in the context of reconstructing scenes from publically available images. This research is motivated by the difficulty in defining building parts as basic units, which are useful for larger areas, the high costs for directly acquiring terrestrial and aerial LiDAR data, and the need to provide textured scenes and hence avoids semantic structural descriptions; for pose estimation techniques for very large number of images see (Snavely et al., 2006; Frahm et al., 2010); for dense surface reconstruction see (Furukawa and Ponce, 2010; Jancosek and Pajdla, 2011; Langguth et al., 2016). 3D surfaces of high fidelity and sufficient density, however, are an ideal basis for deriving semantically rich building descriptions, the topic of the next section.

C. Image Interpretation with Graphical Models

Deriving maps from images (including range images, e.g., LiDAR measurements) – a central task of photogrammetric research – by means of automatic image interpretation techniques still is in a premature state.

We addressed the problem of image interpretation for generating structured scene descriptions using building façades as exemplary domain. Façades show a wide variety in parts (doors, windows, balconies), structure (repetitions, symmetry, alignment) and appearance (local shadows, reflections, vegetation). Due to their mostly two-dimensional character modelling regularities is simpler than when dealing with general 3D building structures. We investigated two approaches: data driven semantic image segmentation using graphical models, especially conditional random fields (CRFs) (Korč, 2012; Yang, 2011), and model driven façade reconstruction using marked point processes (MPPs) (Wenzel, 2016; Wenzel and Förstner, 2016). We only discuss the model-driven approach.

The model driven reconstruction (Wenzel, 2016; Wenzel and Förstner, 2016) starts from rectified images, assuming the scale to be known. The model is a marked point process where façades consist of façade elements (doors, windows, balconies) represented as rectangles. The interpretation uses a reversible jump Markov chain Monte Carlo (rjMCMC) hypothesis and test paradigm. The geometric properties of the elements and their spatial relations are learned from training data. The data term of the energy function depends on a probabilistic object related classification; see Fig. 6. For the bottom example, observe the wrong heights of the windows, the confusion of windows and balconies and the detection of windows, where the ground truth does not indicate them; these errors can be explained by a too weak prior on the neighbourhood relations and the lack of long range interactions between the facade elements. The empirical evaluation of the method leads to confusion tables, which contain estimated conditional probabilities, for which confidence intervals can be given. In order to arrive at reasonable intervals in case the empirical probability is 0 or 1, a weak Dirichlet prior for the multinomial distribution of these empirical probabilities can be used, see Table I.



Fig. 6. Façade reconstruction based on a marked point process for façade elements (image, ground truth, reconstruction), CPU-time appr. one hour. from (Wenzel, 2016)

TABLE I

Top: Confusion matrix for the façade type 'city houses' (with classes, background, window and balcony); **Bottom:** Corresponding 99%-confidence regions in % for the probabilities using a weak Dirichlet prior $\mathcal{D}(\alpha)$ with, e.g., $\alpha = [1, 0.01, 0.01]$ for the first row. This yields more reasonable intervals for cases where the empirical probability is 0 or 1; from (Wenzel, 2016)

				predictio	on		
			bg	win	balc		
	truth	bg	0	5	0		
		win	6	146	0		
		balc	0	3	40		
background		window			balcony		
0.11 < 16.7 < 65	.3	34.5 <	83.1	< 99.9	0.0	0 < 0.17 < 9.5	6
1.03 < 3.93 < 9.0	04	90.9 <	96.1	< 99.0	0.0	0 < 0.01 < 0.3	6
0.00 < 0.02 < 1.2	27	0.81 <	6.84	< 19.8	80.	<i>l</i> < 93.1 < 99.2	!

We observed the typical strengths and weaknesses of data driven and model driven methods. Data driven methods are fast, can adapt locally to the image information and are versatile. This refers not only to locally connected Markov random fields (MRF), which, latest since grab-cut (Rother et al., 2004), pushed research in semantic segmentation, see the review (Zhu et al., 2015). This also holds for (1) fully connected MRFs, e.g., (Krähenbühl and Koltun, 2011; Ristovski et al., 2013; Li and Yang, 2016), which, due to their special assumption on the potentials, easily achieve real time (Cheng et al., 2015) while still being competitive, (2) for autocontext models, which aim at sequentially gathering new context by using features of previous interpretations, see e.g., (Tu, 2008; Jampani et al., 2015; Gadde et al., 2016), but even more also (3) for convolutional neural networks, e.g., (Farabet et al., 2013; Long et al., 2014; Marmanisa et al., 2016). The techniques have also been applied successfully to semantically segmenting point clouds, e.g., (Adan et al., 2011; Tamke et al., 2014; Ochmann et al., 2016). Graphical models may be linked to logical programming via Markov logical networks (Richardson and Domingos, 2006). They allow for a

mixture of crisp and soft formulas Fierens et al. (2014). They are used for event and face recognition in image sequence analysis (Tran and Davis, 2008; Chechetka et al., 2010), for text understanding (Poon, 2011), for the interpretation of images of chemical structures (Frasconi et al., 2014), and for scene interpretation (Xu and Petrou, 2010).

Model driven methods allow to explicitly model long range constraints. This in a first place holds for models based on grammars (Zhu and Mumford, 2006) and marked point processes (Ortner et al., 2007). Grammars are regularly for city modeling (Dick et al., 2002; Talton et al., 2011; Martinović and Gool, 2013; Liu et al., 2014; Schwarz and Müller, 2015) or for roof extraction (Huang et al., 2011). They allow learning, as for indoor scenes (Liu et al., 2014), for façades (Ripperda and Brenner, 2009; Fan et al., 2014; Wu et al., 2014), for building layouts (Bao et al., 2013), or for architectural styles (Talton et al., 2012). In image interpretation marked point processes are used for building extraction (Ortner et al., 2007), for road network extraction (Chai et al., 2013), or more geometric feature extraction (Lafarge et al., 2008). The generality of these models requires costly sampling methods for (approximately) finding optimal interpretations, which, however, allow for parallelization (Wilkinson, 2006).

The *integration of data and model driven methods* has always been the key to successful interpretations. Early approaches, such as (Mohan and Nevatia, 1989), used perceptual grouping techniques for providing candidate regions for object detection, here detecting buildings in aerial images. The same flavour can be found in recent work on simultaneous segmentation and detection (Hariharan et al., 2014), where the region proposals are refined after classification in order to obtain more accurate region boundaries.

A probably first integration of Markov logical networks and stochastical grammars for interpreting façades from point clouds is described in (Dehbi et al., 2016), see Fig. 7. The



Fig. 7. Façade model with stochastical grammar and Markov logical network. **Upper left:** An instance of the grammar. **Upper right:** Some probabilistic rules of the grammar. **Lower left:** Some probabilistic relations of the Markov logic network; from (Dehbi et al., 2016)

partonomy of the façade is represented in stochastic attributed grammatical rules, which capture the geometric properties and relations between the parts, see Fig. 7, upper right. Additional constraints are represented as predicates, which due to the



Fig. 8. Interpretation of point clouds of façades with stochastic grammars and Markov logical networks. **Top row:** Image of façade. **Second row:** Point cloud with holes. **Third row:** Data driven interpretation of point cloud; windows: green), doors: pink, façade: white. **Last row:** Model driven interpretation. Observe the predictive power of the model; from (Dehbi et al., 2016)

diversity of the training data are give a probability. Relations between these predicates establish the Markov logic network, see Fig. 7, lower left. The interpretation of the point cloud starts with detecting basic parts of the façade. Deficiencies such as missing parts, or wrong alignments are then corrected using the prior mode, see Fig. 8.

III. SOME CURRENT PROBLEMS

This section discusses a few problems which regularly appear when developing methods for automatic interpretation of man made scenes. They address the choices we have when modelling the imaging process with the goal to solve the inverse problem, namely to recover scene information from images. Specifically, they refer to the model of the image signal, the relation between discrete and continuous geometry, the integration of crisp and soft prior knowledge, and the type of uncertainty of events and their meaning.

A. Physical and Phenomenological Signal Models

The basic steps for image orientation and building reconstruction, as the examples showed, often use methods for edge and contour detection, which essentially depend on the assumed image model. A classical model for the observed intensities g(i) in an image starts from the photon counts N(i)at each pixel in k channels: the two k-vectors $g(i) \propto N(i)$. This basic assumption leads to several problems, when following classical image processing procedures:

• How to exploit colour theory for non-RGB imagery? Colour theory models are a phenomenological and model visual perception of colours and its peculiarities, such as colour definition or colour constancy, or it models colour printing. For the majority of images available it may be useful: however, the analysis of images with more than three colours, even of hyperspectral images, the basic physical model appears to be the appropriate start. Improvements of classifiers using other than the original RGB signal result from reduced correlations, which are preferred by models which treat features as uncorrelated. Models, which take the – in principle arbitrary – distribution of the three colours into account, would not gain from colour transformations.

- Image intensities, being proportional to photon counts, are positive values. Representing a spatial intensity g(i) as a sum of basis functions which are not non-negative, as when applying Fourier or Wavelet analysis, appears to be physically meaningless. Nevertheless, spectral methods have shown to be very successful.
- Since perception is logarithmic, a simple way out would be to work with the logarithms of the intensities, as proposed by (Koenderink and Doorn, 2002), who motivates it by the logarithmic perception of intensity.² The representation of the positive function then would be similar to the exponential family of densities, see (Borwein and Huang, 1995) and the generalisation in (Fasino, 2002), e.g., when assuming a continuous image domain, $f_l(x) = \log(g(x)) = \int G(u) \exp(2\pi i x u) du$, where G(u) is the Fourier transform of g(x).
- The statistical model of the observed intensities, being proportional to the photon counts, is a Poisson distribution. Then the variance of the intensity increases linearly with the intensity, omitting thermal noise and non-linearities of the sensor. Using a simple box-filter for smoothing implicitly assumes the intensities to have the same variance in the chosen neighbourhood, which does not hold. Checking the gradient magnitude for detecting edges, which is a classification task, should take the variance, i.e., the intensity level into account. Alternatively, the signal could be variance normalized, in the most simple case using a square root point transformation $f_s(x) = \sqrt{g(x)}$, since the normalized signal $f_s(x)$ then has constant noise variance, see (Förstner, 2000; Jähne and Schwarzbauer, 2016). Many algorithms for keypoint detection could gain from such a transformation, leading to less keypoints in bright and more keypoints in dark areas of the image. This type of transformation also is motivated by the sensitivity of visual perception to image coding, see (Mannos and Sakrison, 1974).

It would be desirable to have an integrating model for intensity signals in order to allow for efficient statistical, physical and (spatial) spectral analysis. A scale analysis of the factors resulting from non-negative signal factorization may play a guiding role.

B. Discrete and Continuous Geometry

Recovering man-made objects aims at some geometric description of the object's boundary, which usually is represented as an aggregation of continuous surface regions in 3D. The expected image of such a piecewise surface is a partitioning

 $^{^{2}}$ (Koenderink and Doorn, 2002) in addition allow for affine transformations of the logarithm of the intensity.

of the image region with piecewise boundaries, where not all intensity edges necessarily need to have two distinct regions as neighbours. The observed image grid as observed 3D point clouds are discrete. Hence, the reconstruction of the continuous 3D surface regions and their boundaries consists (1) of the topologically consistent identification of these boundaries and (2) the geometrically consistent determination of the form of the surfaces and the boundaries. An example where boundaries may not be detectable due to lighting conditions and may lead to violations of the image model as is given in Fig. 9.



Fig. 9. Topological relations for polyhedra and their ideal and real (extracted) images. Left: Image of a vertical edge appearing with zero gradient (St. Michaelis church, Hamburg). The following: Example image and entity relation diagram with range of multiplicities. Mid left: Polyhedral boundary; edges (E) may must have two neighbouring regions (R) in 3D. Mid: Ideal image of polyhedral boundary, admitting zero gradient edges; edges may have 1 to 2 neighbouring regions in the ideal image. Right: Real image from partitioning; edges (E) must have two neighbouring regions (R) in the partitioning. In all cases edges have two neighbouring (end) points (P). Obviously, the ideal image does not follow the winged-edge representation

The recovery of a consistent boundary description is underconstrained, unless the point density (of the grid or the point cloud) is sufficiently high and the boundary lines fulfil certain regularities, see e.g., the sampling theorem for recovering region boundaries (Meine et al., 2009), architectural models (Pottmann et al., 2015), not necessarily based on triangular meshes (Liu et al., 2006; Kovacs et al., 2011).

Moreover, grid-based methods, such as Markov random fields, do not allow to include prior knowledge about the straightness of boundaries. Therefore, algorithms for finding consistent polygonal boundaries mostly contain ad hoc rules, cannot include statistical prior information about the observed points, and – due to the occurrence of structural errors – are difficult to be evaluated.

This touches the integration of bottom-up and top-down procedures, discussed above: geometric entities, such as polygons or polygon networks, being mid-level structures, require a statistically coherent modelling of both, their appearance – for bottom-up hypothesis building – as well as their geometric and neighbourhood relations – for top-down prediction; this appears like bi-directional search in the solution space. Reducing the costs for sampling from large energy models, such as with MCMC, is described in (Papandreou and Yuille, 2011).

C. Crisp and Soft Prior Knowledge

Handling both, crisp and soft prior knowledge, as prior is essential (not only) for interpreting images of man-made scenes. Taxonomies and partonomies of objects and spatial relations, such as parallelity, play a central role in semantic modelling, see the discussion of the role of semantics for games in (Tutenel et al., 2008). The uncertainty of observations and models and the success of probabilistic models is ubiquitous.

It is less clear how those parts of the model, which are certain (in a probabilistic sense), are handled in a principled manner: i.e., explicitly. *Geometric relations* in multi-view analysis usually are hard coded; algebraic methods, such as Gröbner bases, though a research topic on its own, increasingly are used to derive solutions, but are not integrated into systems, where the task is not fixed. Attempts to use algebraic methods for more generic tasks, have been intensively discussed in the late eighties, see (Kapur and Mundy, 1989; Mundy et al., 1998). Methods which detect regularities and use them for the enhancement of 2D and 3D objects, such as in (Brenner, 2005; Meidow et al., 2016), have to face the inconsistency of individual hypothesis tests or the explosion of computational complexity – prior to finding bases for the constraints, which then can be applied.

Partonomies and taxonomies are increasingly used for improving categorization (Marszalek and Schmid, 2007; Griffin and Perona, 2008). Following (Zweig and Weinshall, 2007), the simultaneous classification of a category and a subcategory is significantly better than the individual classification. Explicitly classifying image galleries, i.e., ensembles of images, into a given taxonomy (derived from Wikipedia) is adressed by (Kramer et al., 2012). The images in the data base IMAGENET (Deng et al., 2009) are organized in a semantic hierarchy (WordNet), supporting benchmarking of classifiers which can exploit this knowledge. Since ImageNet is based on the ontology of WordNet it would be desirable to have the concepts around 'building' for interpreting outdoor and indoor images included in ImageNet. Since WordNet is focussed on function of notions and does not include any concepts for geometric or material the link between semantic, geometric, and radiometric models still remains to be established, e.g., for the domains 'building' and 'road', possibly exploiting grammars, marked point processes, or Markov logic networks, see the example above.

In this context two questions arise. First, what are the adequate methods to learn the models, i.e., the geometric and semantic relations? Learning the structure and the parameters of probabilistic logic, where clauses are attached with a probability, may be based on measuring the success of data base queries Gutmann et al. (2008); Fierens et al. (2014). Learning structures can use the development in kernel methods, which allow to address all types of structures: multi-label, with taxonomies, label-sequence-learning, sequence of operations alignment, natural language parsing, see (Tsochantaridis et al., 2004).

Second, what are efficient interpretation processes? There exist several methods to derive statistically interpretations based on crisp and uncertain information, e.g., using probabilistic logic programming, statistical relational learning, or Markov logic, see the overview in (Fierens et al., 2014) and the Dagstuhl Seminar on *Logic and Probability for Scene Interpretation*; see (Neumann et al., 2008). Attempts to increase efficiency use a reduced language e.g., (Domingos and Webb, 2012), or apply sampling techniques, e.g., (Poon and Domingos, 2006; Beltagy and Mooney, 2014). Except for a few examples, e.g., (Zhu et al., 2015; Dehbi et al., 2016), see above, the techniques are not yet exploited for analysing images, especially of man-made objects.

D. Uncertainty and Vagueness

Decision making using classifiers always has assumed that data, models, and decisions are uncertain. However, the process and the result of classifiers often do not reflect this uncertainty.

First, many classifiers only report the most likely class for each object in a 'winner takes all' habit. This does not support the need of a user to know the uncertainty of the decision. Even giving a confusion matrix, often is not sufficient, as the estimated conditional probabilities are estimates, and hence are uncertain. Giving confidence regions as in Table I, would be a first remedy. Results of (Roscher, 2012; Roscher et al., 2012) indicate, that import vector machines (Zhu and Hastie, 2001) yield more reliable posterior probabilities than the output of support vector machines, when transforming their output into probabilities (Platt, 1999). Since the output of classifiers often is used for generating the potentials of Markov random fields, their quality may have a decisive impact.

The uncertainty of semantic segmentation cannot be represented with confusion tables, as the space of segmentations is far too large, why indicating the uncertainty of the boundaries appears a reasonable approach, see e.g., (Kendall et al., 2015; Kampffmeyer et al., 2016) both using deep convolutional network.

Second, many classifiers assume that each object belongs to one of the presumed classes, possibly a rejection class. This has been found to be over-simplistic. Images with complex content may belong to different classes, e.g., a natural scene may simultaneously be classified as mountain area and beach area, if ingredients (key features) for both classes can be detected, and the designer of the classifier intends such an overlap of classes, see (Arbelaez et al., 2011; Roth and Fischer, 2007) and the review (Sorower, 2010).

Third, the classes themselves are difficult to separate, e.g., the two classes low vegetation and high vegetation in the 'Large Scale Point Cloud Classification benchmark'.³ This type of uncertainty in the definition of classes was the motivation for developing fuzzy models (Zadeh, 1965, 1975), where each object may belong to a class according to some membership value. The heavy debate on the relation between fuzzy theory and probability theory is reflected and resolved in the key paper by Dubois and Prade (Dubois and Prade, 1993): The semantic distinction between the vagueness/fuzziness of the notion of an event and the uncertainty/likelihood of the existence or appearance of an event indicates the two notions to be orthogonal; integrating both concepts, while keeping their key properties, such as (Zadeh, 1968, 1975; Navara, 2005), still seems to have no canonical solution.

Anyhow, when taking into account the necessity to handle non-unique ground truth in benchmarks (Martin et al., 2001), to deal with occlusions,⁴, and to take vaguely defined classes into account, when evaluating classifiers (see (Everingham et al., 2015)), then the number of papers addressing fuzzy logic on international conferences such as ICPR, ECCV, and ICCV, being below 0.5 % on an average, appears to be very low.

IV. A FUTURE FOR LEARNING BUILDING MODELS

The paper has addressed various aspects of interpreting images of man-made structures, especially of buildings. It focused on methods, which reflect the underlying models of the imaging, analysis and interpretation processes and which hence allow the user of such a system to make decisions, thus to understand the image content in an appropriate manner. The problems, mentioned in the last section, all are caused by insufficiencies or incompatibilities of simultaneously applied models. The tools to solve or overcome these problems appear to be available.

We discussed the dichotomy of discriminative and generative models, both having their advantages. Discriminative models are efficient in obtaining quantitatively good results, while generative models are powerful in elucidating structured semantics. The dichotomy is best seen in semantic segmentation: The partitioning of image into relevant regions requires a process which is at the same time data and model driven. This motivates the structuring of the interpretation/understanding task as in Fig. 10.



Fig. 10. Metamodel for interpretion and understanding of image data (left ellipse). Preknowledge (right ellipse), e.g., in the form of grammars or marked point processes (MPP), are necessary, in order to capture the envisaged meaning of the interpretation. Intermediate structures (mid ellipse) may be represented and analysed e.g., by Markov random fields (MRF), conditional random fields (CRF) or Markov logic networks (MLN). These structures are *simultaneously* predicted from the preknowledge and from the data by – possibly semantic – segmentation. The parameters of all processes (indicated with upright letters and arrows) are trained using techniques from machine learning (ML). The control of the complete process is an open problem

All processes can gain from the interaction of recognition, reconstruction and re-organization, proposed in (Malik et al., 2015). Generative models also are directly amenable to incremental learning. On the other hand the speed of current neural network classification and regression tools, which does

³See http://www.semantic3d.net

⁴See e.g., the annotation rules in the PASCAL http://host.robots.ox.ac.uk/ pascal/VOC/voc2008/guidelines.html.

in no way correspond to the generally long training times, contrasts to the fast training times of explicit semantic models, such as grammars of marked point processes, which are often much slower in reasoning. Attempts to use neural network priors for one-shot learning, such as (Fu et al., 2015), are promising. The flexibility of multi-layer neural networks also needs to be compared with the rich representation of the scattering transform (Bruna and Mallat, 2012; Mallat, 2016), which code higher moments of the underlying signal, and have been applied in face recognition (Chang et al., 2012), used for graphs (Chen et al., 2014), and enriched by rotation invariant kernels (Tolias et al., 2015).

The trend to have very large and rich bodies of image data for benchmarking can be interpreted as extensionally defining what an image is, instead of intentionally modelling images by power spectra, higher order characteristics, or random fields.

A future for learning highly structured models may be based on the available basic technology which not yet is exploited for establishing rich geometric and semantic building models.

REFERENCES

- Adan, A., X. Xiong, B. Akinci, and D. Huber (2011, June). Automatic Creation of Semantically Rich 3D Building Models from Laser Scanner Data. In Proceedings of the International Symposium on Automation and Robotics in Construction (ISARC).
- Arbelaez, P., M. Maire, C. Fowlkes, and J. Malik (2011, May). Contour Detection and Hierarchical Image Segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.* 33(5), 898–916.
- Bao, F., D.-M. Yan, N. J. Mitra, and P. Wonka (2013). Generating and Exploring Good Building Layouts. ACM Transactions on Graphics 32.4, 122.
- Beltagy, I. and R. J. Mooney (2014). Efficient Markov Logic Inference for Natural Language Semantics. In Workshop at the Twenty-Eighth AAAI Conference on Artificial Intelligence.
- Borwein, J. M. and W. Z. Huang (1995). A Fast Heuristic Method for Polynomial Moment Problems with Boltzmann-Shannon Entropy. *SIAM Journal on Optimization* 5(1), 68–99.
- Brachmann, E., F. Michel, A. Krull, M. Y. Yang, S. Gumhold, and C. Rother (2016). Uncertainty-Driven 6D Pose Estimation of Objects and Scenes from a Single RGB Image. In Proc. of Conf. on Computer Vision and Pattern Recognition.
- Braun, C., T. H. Kolbe, F. Lang, W. Schickler, V. Steinhage, A. Cremers, W. Förstner, and L. Plümer (1995). Models for Photogrammetric Building Reconstruction. *Computer & Graphics 19*(1), 109–118.
- Brenner, C. (2000). Towards Fully Automatic Generation of City Models. In *International Archives of Photogrammetry and Remote Sensing*, Volume XXXIII, pp. 85–92.
- Brenner, C. (2005). Constraints for Modelling Complex Objects. In Int. Archives for Photogrammetry and Remote Sensing, Volume XXXVI, 3/W24.
- Brooks, R. (1983). Model-based 3-D interpretation of 2-D images. *IEEE T-PAMI* 5(2), 140–150.
- Brooks, R. A., R. Creiner, and T. O. Binford (1979). The ACRONYM Model-based Vision System. In *Proceedings of the 6th International Joint Conference on Artificial Intelligence - Volume 1*, IJCAI'79, San Francisco, CA, USA, pp. 105–113. Morgan Kaufmann Publishers Inc.
- Bruna, J. and S. Mallat (2012). Invariant Scattering Convolution Networks. CoRR abs/1203.1513.
- Chai, D., W. Förstner, and F. Lafarge (2013). Recovering Line-Networks in Images by Junction-Point Processes. In *Conf. on Computer Vision and Pattern Recognition.*

- Chang, K.-Y., C.-F. Lin, C.-S. Chen, and Y.-P. Hung (2012). Applying scattering operators for face recognition: A comparative study. In *Pattern Recognition (ICPR), 21st International Conference on.*
- Chechetka, A., D. Dash, and M. Philipose (2010). Relational Learning for Collective Classification of Entities in Images. In Proceedings of the 6th AAAI Conference on Statistical Relational Artificial Intelligence, AAAIWS'10-06, pp. 7–12. AAAI Press.
- Chen, X., X. Cheng, and S. Mallat (2014). Unsupervised Deep Haar Scattering on Graphs. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K. Weinberger (Eds.), Advances in Neural Information Processing Systems 27, pp. 1709–1717. Curran Associates, Inc.
- Cheng, M.-M., V. A. Prisacariu, S. Zheng, P. H. S. Torr, and C. Rother (2015). DenseCut: Densely Connected CRFs for Realtime Grab-Cut. *Computer Graphics Forum* 34(7), –.
- Dehbi, Y., F. Hadiji, G. Gröger, K. Kersting, and L. Plümer (2016). Statistical Relational Learning of Grammar Rules for 3D Building Reconstruction. *Transactions in GIS*.
- Deng, J., W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei (2009). ImageNet: A Large-Scale Hierarchical Image Database. In Proc. of Conf. on Computer Vision and Pattern Recognition.
- Dick, A. R., P. H. S. Torr, and R. Cipolla (2002). A Bayesian Estimation of Building Shape using MCMC. In *Computer Vision– ECCV 2002*, Volume 2351 of *LNCS*, pp. 574–575. Springer.
- Domingos, P. and W. A. Webb (2012). A tractable first-order probabilistic logic. In *Proceedings of the Twenty-Sixth AAAI Conference on Artificial Intelligence*, AAAI'12, pp. 1902–1909. AAAI Press.
- Dubois, D. and H. Prade (1993). Fuzzy sets and probability: misunderstandings, bridges and gaps. In *International Conference* on Fuzzy Systems, Volume 2, pp. 1059–1068.
- Everingham, M., S. M. A. Eslami, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman (2015, January). The Pascal Visual Object Classes Challenge: A Retrospective. *International Journal* of Computer Vision 111(1), 98–136.
- Fan, L., P. Musialski, L. Liu, and P. Wonka (2014, November). Structure Completion for Facade Layouts. ACM Trans. Graph. 33(6), 210:1–210:11.
- Farabet, C., C. Couprie, L. Najman, and Y. LeCun (2013, August). Learning Hierarchical Features for Scene Labeling. *IEEE Trans*actions on Pattern Analysis and Machine Intelligence.
- Fasino, D. (2002, March). Approximation of nonnegative functions by means of exponentiated trigonometric polynomials. J. Comput. Appl. Math. 140(1-2), 315–329.
- Fidler, S. and R. Urtasun (2015). 3D Indoor Scene reconstruction. http://www.cs.toronto.edu/~fidler/slides/CVPR15tutorial/.
- Fierens, D., G. Van Den Broeck, J. Renkens, D. Shterionov, B. Gutmann, I. Thon, G. Janssens, and L. de Raedt (2014). Inference and Learning in Probabilistic Logic Programs using Weighted Boolean Formulas. *Theory and Practice of Logic Programming* 15(3), 358– 401.
- Fischer, A. (2000). Automatische Gebäuderekonstruktion mittels parametrisierter Komponenten. Ph. D. thesis, Institut für Informatik, Universität Bonn, eingereicht.
- Fischer, A., T. H. Kolbe, F. Lang, A. B. Cremers, W. Förstner, L. Plümer, and V. Steinhage (1998). Extracting Buildings from Aerial Images Using Hierarchical Aggregation in 2D and 3D. *CVIU* 72(2), 185–203.
- Förstner, W. (1988). Model Based Detection and Location of Houses as Topographic Control Points in Digital Images. In *Intl. Archives* of Photogrammetry and Remote Sensing, Volume 27. XVIth ISPRS Congress, Kyoto.
- Förstner, W. (2000). Image Preprocessing for Feature Extraction in Digital Intensity, Color and Range Images. In *Geomatic Methods* for the Analysis of Data in Earth Sciences, Volume 95/2000 of Lecture Notes in Earth Sciences, pp. 165–189. Springer.
- Frahm, J.-M., P. Fite-Georgel, D. Gallup, T. Johnson, R. Raguram,

C. Wu, Y.-H. Jen, E. Dunn, B. Clipp, S. Lazebnik, and M. Pollefeys (2010). Building Rome on a Cloudless Day. In *Proceedings of the 11th European Conference on Computer Vision: Part IV*, ECCV'10, Berlin, Heidelberg, pp. 368–381. Springer-Verlag.

- Frasconi, P., F. Gabbrielli, M. Lippi, and S. Marinai (2014). Markov Logic Networks for Optical Chemical Structure Recognition. J. of Chemical Information and Modeling 54(8), 2380–2390.
- Fu, J., S. Levine, and P. Abbeel (2015). One-Shot Learning of Manipulation Skills with Online Dynamics Adaptation and Neural Network Priors. *CoRR abs/1509.06841*.
- Fua, P. and A. J. Hanson (1987). Reseguentation Using Generic Shape Locating General Cultural Objects. *Pattern Recognition Letters* 5, 243–252.
- Fuchs, C. and W. Förstner (1995). Polymorphic Grouping for Image Segmentation. In Proc. 5th ICCV, pp. 175–182. IEEE Computer Society Press.
- Furukawa, Y. and J. Ponce (2010, August). Accurate, Dense, and Robust Multiview Stereopsis. *IEEE Trans. Pattern Anal. Mach. Intell.* 32(8), 1362–1376.
- Gadde, R., V. Jampani, R. Marlet, and P. V. Gehler (2016). Efficient 2D and 3D Facade Segmentation using Auto-Context. *CoRR abs/1606.06437*.
- Griffin, G. and P. Perona (2008). Learning and using taxonomies for fast visual categorization. In *Proc. of Conf. on Computer Vision and Pattern Recognition*.
- Grimson, W. E. L. and T. Lozano-Perez (1987). Localizing Overlapping Parts by Searching the Interpretation Tree. *IEEE T-PAMI 9*(4), 469–482.
- Gutmann, B., A. Kimmig, K. Kersting, and L. d. Raedt (2008). Parameter learning on probabilistic databases – A least squares approach. In *European Conference on Machine Learning*.
- Haala, N. and G. Vosselman (1992). Recognition of Road and River Patterns by Relational Matching. In *Intl. Archives for Photogrammetry*, Volume XXIX, Part B3, pp. 969–975.
- Haralick, R. M. and L. G. Shapiro (1979). The Consistent Labeling Problem: Part 1. *IEEE T-PAMI 1*(2), 173–184.
- Haralick, R. M. and L. G. Shapiro (1980). The Consistent Labeling Problem: Part 2. *IEEE T-PAMI* 2(3), 193–203.
- Hariharan, B., P. Arbelaez, R. B. Girshick, and J. Malik (2014). Simultaneous Detection and Segmentation. *CoRR abs/1407.1808*.
- Herman, M. and T. Kanade (1986). Incremental Reconstruction of 3D Scenes from Multiple, Complex Images. *Artificial Intelligence 30*, 289–341.
- Heuel, S. and T. H. Kolbe (2001, July). Building Reconstruction: The Dilemma of Generic Versus Specific Models. *Künstliche Intelligenz 3*, 57–62.
- Huang, H., C. Brenner, and M. Sester (2011). 3D Building Roof Reconstruction from Point Clouds via Generative Models. In Proceedings of the 19th ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems, GIS '11, New York, NY, USA, pp. 16–24. ACM.
- Huertas, A. and R. Nevatia (1988). Detecting Buildings in Aerial Images. *Computer Vision, Graphics, and Image Processing* (41), 131–152.
- Jampani, V., R. Gadde, and P. V. Gehler (2015). Efficient facade segmentation using auto-context. In 2015 IEEE Winter Conference on Applications of Computer Vision, WACV 2015, Waikoloa, HI, USA, January 5-9, 2015, pp. 1038–1045.
- Jancosek, M. and T. Pajdla (2011). Multi-view Reconstruction Preserving Weakly-supported Surfaces. In Proc. of Conf. on Computer Vision and Pattern Recognition, CVPR '11, Washington, DC, USA, pp. 3121–3128. IEEE Computer Society.
- Jähne, B. and M. Schwarzbauer (2016). Noise equalisation and quasi loss-less image data compression – or how many bits needs an image sensor? *Technisches Messen, De Gruyter 83*, 16–24.
- Kampffmeyer, M., A.-B. Salberg, and R. Jenssen (2016). Semantic Segmentation of Small Objects and Modeling of Uncertainty in

Urban Remote Sensing Images Using Deep Convolutional Neural Networks. In CVPR Workshop 'Visual Analysis of Satellite to Street Imagery Workshop'.

- Kapur, D. and J. L. Mundy (1989). Geometric Reasoning and Artificial Intelligence: Introduction to the Special Volume. In *Geometric Reasoning*, pp. 1–14.
- Kar, A., S. Tulsiani, J. Carreira, and J. Malik (2014). Category-Specific Object Reconstruction from a Single Image. *CoRR abs/1411.6069*.
- Kendall, A., V. Badrinarayanan, and R. Cipolla (2015). Bayesian segnet: Model uncertainty in deep convolutional encoder-decoder architectures for scene understanding. *CoRR abs/1511.02680*.
- Koenderink, J. J. and A. J. v. Doorn (2002). Image Processing Done Right. In *Proceedings of the 7th European Conference on Computer Vision-Part I*, ECCV '02, London, UK, UK, pp. 158– 172. Springer-Verlag.
- Kolbe, T. H. (2000). Identifikation und Rekonstruktion von Gebäuden in Luftbildern mittels unscharfer Constraints. Shaker Verlag: Shaker.
- Kolbe, T. H., L. Plümer, and A. B. Cremers (2000, January). Identifying Buildings in Aerial Images Using Constraint Relaxation and Variable Elimination. *IEEE Intelligent Systems* 15(1), 33–39.
- Korč, F. (2012). Tractable Learning for a Class of Global Discriminative Models for Context Sensitive Image Interpretation. Ph. D. thesis, Department of Photogrammetry, University of Bonn, http://hss.ulb.uni-bonn.de/2012/3010/3010.htm.
- Kovacs, D., A. Myles, and D. Zorin (2011). Anisotropic quadrangulation. *Computer Aided Geometric Design* 28(8), 449 – 462.
- Krähenbühl, P. and V. Koltun (2011). Efficient Inference in Fully Connected CRFs with Gaussian Edge Potentials. In J. Shawe-Taylor, R. S. Zemel, P. L. Bartlett, F. Pereira, and K. Q. Weinberger (Eds.), Advances in Neural Information Processing Systems 24, pp. 109–117. Curran Associates, Inc.
- Kramer, G., G. Bouma, D. Hendriksen, and M. Homminga (2012). Classifying Image Galleries into a Taxonomy Using Metadata and Wikipedia. In Proceedings of the 17th International Conference on Applications of Natural Language Processing and Information Systems, NLDB'12, Berlin, Heidelberg, pp. 191–196. Springer-Verlag.
- Lafarge, F., G. Gimel'farb, and X. Descombes (2008). Geometric Feature Extraction by a Multi-Marked Point Process. *Transaction* on Pattern Analysis and Machine Intelligence 32, 1597–1609.
- Lake, B. M., T. D. Ullman, J. B. Tenenbaum, and S. J. Gershman (2016). Building Machines That Learn and Think Like People. *CoRR abs/1604.00289*.
- Lang, F. and W. Förstner (1996). Surface Reconstruction of Man-Made Objects using Polymorphic Mid-Level Features and Generic Scene Knowledge. Zeitschrift für Photogrammetrie und Fernerkundung 6, 193–201.
- Langguth, F., K. Sunkavalli, and M. G. Sunil Hada and (2016). Shading-aware Multi-view Stereo. In Proc. of European Conference on Computer Vision.
- Li, W. and M. Y. Yang (2016). Efficient semantic segmentation of man-made scenes using fully-connected conditional random fields. In Archives of ISPRS.
- Liu, T., S. Chaudhuri, V. G. Kim, Q. Huang, N. J. Mitra, and T. Funkhouser (2014, November). Creating Consistent Scene Graphs Using a Probabilistic Grammar. ACM Trans. Graph. 33(6), 211:1–211:12.
- Liu, Y., H. Pottmann, J. Wallner, Y.-L. Yang, and W. Wang (2006, July). Geometric Modeling with Conical Meshes and Developable Surfaces. ACM Trans. Graph. 25(3), 681–689.
- Long, J., E. Shelhamer, and T. Darrell (2014). Fully Convolutional Networks for Semantic Segmentation. CoRR abs/1411.4038.
- Lowe, D. G. (1987). Three-Dimensional Object Recognition from Single Two-Dimensional Images. Artificial Intelligence 31, 355– 395.

- Malik, J., P. Arbelez, J. Carreira, K. Fragkiadaki, R. Girshick, G. Gkioxari, S. Gupta, B. Hariharan, A. Kar, and S. Tulsiani (2015). The Three R's of Computer Vision: Recognition, Reconstruction and Reorganization. *Pattern Recognition Letters* 72, 4–14.
- Mallat, S. (2016). Understanding deep convolutional networks. Philosophical Transactions of the Royal Society of London A: Mathematical, Physical and Engineering Sciences 374(2065).
- Mannos, J. L. and D. J. Sakrison (1974, July). The Effects of a Visual Fidelity Criterion on the Encoding of Images. *IEEE Transactions* on Information Theory IT-20(4), 525–536.
- Marmanisa, D., J. D. Wegner, K. S. S. Gallian and, M. Datcu, and U. Stilla (2016). Semantic segmentation of aerial images with an essemble of CNNs. In *Archives of ISPRS*.
- Marszalek, M. and C. Schmid (2007). Semantic Hierarchies for Visual Object Recognition. In *Proc. of Conf. on Computer Vision and Pattern Recognition*.
- Martin, D., C. Fowlkes, D. Tal, and J. Malik (2001, July). A Database of Human Segmented Natural Images and its Application to Evaluating Segmentation Algorithms and Measuring Ecological Statistics. In *Proc. 8th Int'l Conf. Computer Vision*, Volume 2, pp. 416–423.
- Martinović, A. and L. V. Gool (2013). Bayesian Grammar Learning for Inverse Procedural Modeling. In Proc. of Conf. on Computer Vision and Pattern Recognition.
- Mason, S., M. Baltsavias, and D. Stallmann (1994). High Precision Photogrammetric Data Set for Building Reconstruction and Terrain Modelling. ETH Zürich.
- Meidow, J., H. Hammer, M. Pohl, and D. Bulatov (2016). Enhancement of Generic Building Models by Recognition and Enforcement of Geometric Constraints. In *ISPRS Annals of the Photogrammetry*, *Remote Sensing and Spatial Information Sciences*, Volume III-3.
- Meine, H., U. Köthe, and P. Stelldinger (2009, February). A Topological Sampling Theorem for Robust Boundary Reconstruction and Image Segmentation. *Discrete Appl. Math.* 157(3), 524–541.
- Mohan, R. and R. Nevatia (1989). Using Perceptual Organization to extract 3D-structures. *IEEE T-PAMI 11*, 1121–1139.
- Mottaghi, R., Y. Xiang, and S. Savarese (2015). A Coarse-to-Fine Model for 3D Pose Estimation and Sub-category Recognition. *CoRR abs/1504.02764*.
- Mundy, J. L., O. Faugeras, T. Kanade, C. d'Souza, and M. Sabin (1998). Object Recognition Based on Geometry: Progress over Three Decades. *Philosophical Transactions: Mathematical, Physical and Engineering Sciences* 356, 1213–1231.
- Musialski, P., P. Wonka, D. G. Aliaga, M. Wimmer, L. van Gool, and W. Purgathofer (2013, September). A Survey of Urban Reconstruction. *Computer Graphics Forum* 32(6), 146–177.
- Navara, M. (2005). Probability theory of fuzzy events. In *European* Society for Fuzzy Logic and Technology.
- Neumann, B., A. C. Cohn, D. C. Hogg, and R. Möller (2008). 08091 Abstracts Collection – Logic and Probability for Scene Interpretation. In A. G. Cohn, D. C. Hogg, R. Möller, and B. Neumann (Eds.), *Logic and Probability for Scene Interpretation*, Number 08091 in Dagstuhl Seminar Proceedings, Dagstuhl, Germany. Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik, Germany.
- Ochmann, S., R. Vock, R. Wessel, and R. Klein (2016, February). Automatic Reconstruction of Parametric Building Models from Indoor Point Clouds. *Computers & Graphics* 54, 94–103. Special Issue on CAD/Graphics 2015.
- Oesau, S., F. Lafarge, and P. Alliez (2013, May). Indoor Scene Reconstruction using Primitive-driven Space Partitioning and Graphcut. In *Eurographics Workshop on Urban Data Modelling and Visualisation*, Girona, Spain.
- Ortner, M., X. Descombes, and J. Zerubia (2007, April). Building Outline Extraction from Digital Elevation Models Using Marked Point Processes. *Int. J. Comput. Vision* 72(2), 107–132.

- Papandreou, G. and A. Yuille (2011, November). Perturb-and-map random fields: Using discrete optimization to learn and sample from energy models. In *Proc. IEEE Int. Conf. on Computer Vision* (*ICCV*), Barcelona, Spain, pp. 193–200.
- Platt, J. C. (1999). Probabilistic Outputs for Support Vector Machines and Comparisons to Regularized Likelihood Methods. In Advances in Large Margin Classifiers, pp. 61–74. MIT Press.
- Poon, H. (2011). Markov Logic for Machine Reading. Ph. D. thesis, University of Washington.
- Poon, H. and P. Domingos (2006). Sound and Efficient Inference with Probabilistic and Deterministic Dependencies. In *Proceedings of* the 21st National Conference on Artificial Intelligence - Volume 1, AAAI'06, pp. 458–463. AAAI Press.
- Pottmann, H., M. Eigensatz, A. Vaxman, and J. Wallner (2015). Architectural geometry. *Computers and Graphics* 47, 145–164.
- Richardson, M. and P. Domingos (2006). Markov logic networks. Machine Learning 62(1), 107–136.
- Ripperda, N. and C. Brenner (2009). Evaluation of Structure Recognition using Labelled Facade Images. In J. Denzler, G. Notni, and H. Süße (Eds.), *Pattern Recognition: 31st DAGM Symposium*, pp. 532–541. Springer.
- Ristovski, K., V. Radosavljevic, S. Vucetic, and Z. Obradovic (2013). Continuous Conditional Random Fields for Efficient Regression in Large Fully Connected Graphs. In Proceedings of the Twenty-Seventh AAAI Conference on Artificial Intelligence.
- Roscher, R. (2012). Sequential Learning using Incremental Import Vector Machines for Semantic Segmentation. Ph. D. thesis, University Bonn, http://hss.ulb.uni-bonn.de/2012/3009/3009.htm.
- Roscher, R., W. Förstner, and B. Waske (2012). I2VM: Incremental import vector machines. *Image and Vision Computing 30*(4-5), 263–278.
- Roth, V. and B. Fischer (2007). Improved functional prediction of proteins by learning kernel combinations in multilabel settings. *BMC Bioinformatics* 8, S12 – S12.
- Rother, C., V. Kolmogorov, and A. Blake (2004). GrabCut: Interactive Foreground Extraction using Iterated Graph Cuts. *ACM Transactions on Graphics 23*, 309–314.
- Schickler, W. (1992). Feature Matching for Outer Orientation of Single Images Using 3-D Wireframe Controlpoints. In Intl. Archives of Photogrammetry and Remote Sensing, Volume 29, pp. 591–598. Proc. XVIIth ISPRS Congress, Washington, D. C.
- Schwarz, M. and P. Müller (2015, July). Advanced Procedural Modeling of Architecture. ACM Trans. Graph. 34(4), 107:1– 107:12.
- Shapiro, L. G. and R. M. Haralick (1987, may). Relational Matching. Applied Optics 26(10), 1845–1851.
- Snavely, N., S. M. Seitz, and R. Szeliski (2006, July). Photo Tourism: Exploring Photo Collections in 3D. ACM Trans. Graph. 25(3), 835–846.
- Sorower, M. S. (2010). A Literature Survey on Algorithms for Multilabel Learning. Technical report, Corvallis, OR, Oregon State University.
- Talton, J. O., Y. Lou, S. Lesser, J. Duke, R. Mčch, and V. Koltun (2011, April). Metropolis Procedural Modeling. ACM Trans. Graph. 30(2), 11:1–11:14.
- Talton, J. O., L. Yang, R. Kumar, M. Lim, N. D. Goodman, and R. Měch (2012). Learning Design Patterns with Bayesian Grammar Induction. In Proceedings of the 25th annual ACM symposium on user interface software and technology.
- Tamke, M., I. Blümel, S. Ochmann, R. Vock, and R. Wessel (2014). From Point Clouds to Definitions of Architectural Space. In Fusion - 32nd International Conference on Education and research in Computer Aided Architectural Design in Europe, Volume 2 of eCAADe: Conferences, Newcastle, UK, pp. 557–566. Northumbria University: Northumbria University.
- Tolias, G., A. Bursuc, T. Furon, and H. Jégou (2015, June). Rotation and translation covariant match kernels for image retrieval.

Computer Vision and Image Understanding.

- Tran, S. D. and L. S. Davis (2008). Event modeling and recognition using markov logic networks. In *Proceedings of the 10th European Conference on Computer Vision: Part II*, ECCV '08, Berlin, Heidelberg, pp. 610–623. Springer-Verlag.
- Tsochantaridis, I., T. Hofmann, T. Joachims, and Y. Altun (2004). Support Vector Machine Learning for Interdependent and Structured Output Spaces. In *Proceedings of the Twenty-first International Conference on Machine Learning*, ICML '04, New York, NY, USA, pp. 104–111. ACM.
- Tu, Z. (2008). Auto-context and Its Application to High-level Vision Tasks. In Proc. of Conf. on Computer Vision and Pattern Recognition.
- Tulsiani, S. and J. Malik (2015). Viewpoints and Keypoints. In Proc. of Conf. on Computer Vision and Pattern Recognition, pp. 1510–1519. IEEE.
- Tutenel, T., R. Bidarra, R. M. Smelik, and K. J. D. Kraker (2008, December). The Role of Semantics in Games and Simulations. *Comput. Entertain.* 6(4), 57:1–57:35.
- Vosselman, G. (1992). *Relational Matching*, Volume 628 of *Lecture Notes on Computer Science*. Springer Verlag.
- Vosselman, G. and N. Haala (1992). Erkennung topographischer Paßpunkte durch relationale Zuordnung. Zeitschrift für Photogrammetrie und Fernerkundung 6, 170–176.
- Wenzel, S. (2016). High-Level Facade Image Interpretation using Marked Point Processes. Ph. D. thesis, Department of Photogrammetry, University of Bonn.
- Wenzel, S. and W. Förstner (2016). Facade Interpretation Using a Marked Point Process. In Int. Ann. Photogramm. Remote Sens. (ISPRS'16).
- Wilkinson, D. J. (2006). Parallel bayesian computation, Chapter 18, pp. 481–512. Marcel Dekker/CRC Press.
- Wu, F., D.-M. Yan, W. Dong, X. Zhang, and P. Wonka (2014, July). Inverse Procedural Modeling of Facade Layouts. ACM Trans. Graph. 33(4), 121:1–121:10.
- Xiang, Y., R. Mottaghi, and S. Savarese (2014). Beyond PASCAL: A Benchmark for 3D Object Detection in the Wild. In *IEEE Winter Conference on Applications of Computer Vision (WACV)*.
- Xu, M. and M. Petrou (2010). Learning Logic Rules for Scene Interpretation Based on Markov Logic Networks. In *Proceedings* of the 9th Asian Conference on Computer Vision - Volume Part III, ACCV'09, Berlin, Heidelberg, pp. 341–350. Springer-Verlag.
- Yang, M. Y. (2011). Hierarchical and Spatial Structures for Interpreting Images of Man-made Scenes Using Graphical Models. Ph. D. thesis, Institute of Photogrammetry, University of Bonn, http://hss.ulb.uni-bonn.de/2012/2765/2765.htm.
- Zadeh, L. A. (1965). Fuzzy Sets. Information and Control 8, 338–353.
- Zadeh, L. A. (1968). Probability Measures of Fuzzy Events. J. of Math. Analysis and Applications 23(2), 421–427.
- Zadeh, L. A. (1975). The Concept of a Linguistic Variable and its Application to Approximate Reasoning - I. *Information Sciences* 8, 199–249.
- Zhou, Q.-Y. and U. Neumann (2010). 2.5D Dual Contouring: A Robust Approach to Creating Building Models from Aerial LiDAR Point Clouds. In Proceedings of the 11th European Conference on Computer Vision Conference on Computer Vision: Part III, ECCV'10, Berlin, Heidelberg, pp. 115–128. Springer-Verlag.
- Zhu, H., F. Meng, J. Cai, and S. Lu (2015). Beyond Pixels: A Comprehensive Survey from Bottom-up to Semantic Image Segmentation and Cosegmentation. *CoRR abs/1502.00717*.
- Zhu, J. and T. Hastie (2001). Kernel Logistic Regression and the Import Vector Machine. In *Journal of Computational and Graphical Statistics*, pp. 1081–1088. MIT Press.
- Zhu, S.-C. and D. Mumford (2006, January). A Stochastic Grammar of Images. *Found. Trends. Comput. Graph. Vis.* 2(4), 259–362.
- Zhu, Y., C. Zhang, C. Ré, and L. Fei-Fei (2015). Building a Large-

scale Multimodal Knowledge Base for Visual Question Answering. *CoRR abs/1507.05670.*

Zweig, A. and D. Weinshall (2007). Exploiting Object Hierarchy: Combining Models from Different Category Levels. In *Proc. of Int. Conf. on Computer Vision*.